

AD-A054 555

WISCONSIN UNIV-MADISON MATHEMATICS RESEARCH CENTER

F/G 12/1

SOME OPTIMAL PROPERTIES AND INTERPRETATIONS OF PRINCIPAL COMPO--ETC(U)

MAR 78 R HUDLET, R A JOHNSON

DAA629-75-C-0024

UNCLASSIFIED

MRC-TSR-1836

NL

1 OF 1
AD
A054555



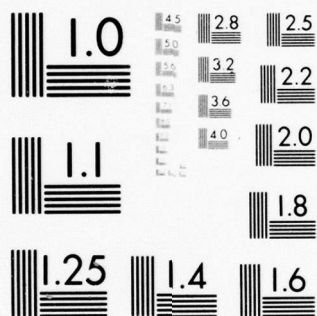
END

DATE

FILMED

7 -78

DDC



MICROCOPY RESOLUTION TEST CHART
NATIONAL BUREAU OF STANDARDS-1963-A

FOR FURTHER TRAN

(13) SC

AD A 054555

MRC Technical Summary Report #1836

SOME OPTIMAL PROPERTIES AND INTERPRETATIONS OF PRINCIPAL COMPONENTS

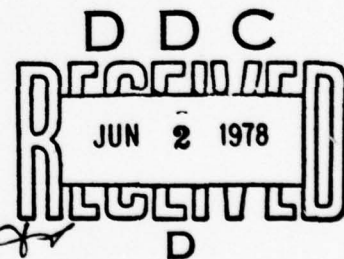
Raul Hudlet and Richard A. Johnson

Mathematics Research Center
University of Wisconsin-Madison
610 Walnut Street
Madison, Wisconsin 53706

March 1978

(Received December 21, 1977)

DDC FILE COPY



Approved for public release
Distribution unlimited

Sponsored by

U. S. Army Research Office
P. O. Box 12211
Research Triangle Park
North Carolina 27709

National Science Foundation
Washington, DC 20550

UNIVERSITY OF WISCONSIN - MADISON
MATHEMATICS RESEARCH CENTER

SOME OPTIMAL PROPERTIES AND INTERPRETATIONS OF PRINCIPAL COMPONENTS

Raul Hudlet and Richard A. Johnson

Technical Summary Report # 1836
March 1978

ABSTRACT

Optimal properties are derived and some new geometrical interpretations given for principal components. Typically, our main results concern the simultaneous minimization of eigenvalues of certain covariance matrices which measure the goodness of an approximation. Many popular criteria like total variance and generalized variance, which are increasing functions of the eigenvalues, are then minimized by the best approximator.

In other situations, the criterion may not be a monotone function of the eigenvalues. In Theorem 3.2, we derive a general optimal class based on the non-negative definite ordering of covariance matrices. Theorem 4.1 gives a result for the sequential selection of principal components. In the final section, we give a new geometrical interpretation of the sample principal components.

AMS (MOS) Subject Classification: 62H25

Key Words: Principal components, Statistical approximations

Work Unit Number 4 (Probability, Statistics, and Combinatorics)

| | |
|--------------------------------------|---|
| ACCESSION FOR | |
| DTIC | White Section <input checked="" type="checkbox"/> |
| DDI | Grey Section <input type="checkbox"/> |
| UNANNOUNCED <input type="checkbox"/> | |
| JUSTIFICATION | |
| BY..... | |
| DISTRIBUTION/AVAILABILITY CODES | |
| Dist. AVAIL. and/or SPECIAL | |
| A | |

SIGNIFICANCE AND EXPLANATION

When a large number of characteristics are measured on a large number of population units, the sheer volume of data can cause problems. It is natural, in such cases, to look for ways to reduce that data to a more manageable form.

One way of doing this is the Principal Component Method, wherein the original problem is replaced by an approximation of far lower dimension. The variables in the approximate problem are certain linear combinations, or principal components, of the variables in the original problem.

Of course, some choices of the exact linear combinations to be used (approximators) will yield more meaningful results than others. Ideally, we would like to retain as much information as possible, and we seek a set of principal components which is optimal from this point of view.

In this paper, we derive some new properties of optimal principal component approximators and obtain some further geometrical insights concerning this method. We also help clarify a potential weakness in this method by determining an even more general class of approximators, that contains those selected by the principal component, and exhibiting a situation where the approximator that would be selected by that method is not optimal.

The responsibility for the wording and views expressed in this descriptive summary lies with MRC, and not with the authors of this report.

SOME OPTIMAL PROPERTIES AND INTERPRETATIONS OF PRINCIPAL COMPONENTS

Raul Hudlet and Richard A. Johnson

1. Review of Previous Work

Let $\underline{X} = (X_1, \dots, X_p)'$ be a random vector with zero expectation and covariance matrix Σ and let $\lambda_1 \geq \dots \geq \lambda_p$ denote the eigenvalues of Σ and $\underline{P}_1, \dots, \underline{P}_p$ a corresponding set of orthonormal eigenvectors. Notice that $\underline{P}_1, \dots, \underline{P}_p$ are uniquely determined (up to multiplication by ± 1) only if $\lambda_1 > \dots > \lambda_p$.

Even though Pearson had encountered principal components as early as 1901, the concept is generally attributed to Hotelling (1933) who was the first to introduce it in a probabilistic framework. Since then Girshick (1936), Anderson (1958), Rao (1964), Darroch (1965), Okamoto and Kanazawa (1968) among others have characterized principal components by different sets of optimality properties.

Girshick (1936) showed that if the components of \underline{X} have variance one, then $\underline{P}_1' \underline{X}$, a first principal component (P.C.), maximizes the sum of the squares of the correlation between $\underline{P}_1' \underline{X}$ and each variate X_i over all possible linear functions $\underline{P}_1' \underline{X}$.

Anderson (1958) established that among the class of linear functions $\underline{P}_1' \underline{X}$ with $\underline{P}_1' \underline{P}_1 = 1$; a first P.C., $\underline{P}_1' \underline{X}$ has maximum variance; $\underline{P}_2' \underline{X}$, a second P.C., has maximum variance among the elements in the class uncorrelated with $\underline{P}_1' \underline{X}$ and so on.

Rao (1964) characterizes the first $k (\leq p)$ principal components as a linear form $\underline{Y} = T \underline{X}$, where T is a $k \times p$ matrix, which minimizes the trace or the Euclidean norm ($\|\underline{M}\|^2 = \sum_{i,j} M_{ij}^2$) of the covariance matrix of the residual of \underline{X} minus the best linear predictor based on \underline{Y} .

Darroch (1965) was the first to characterize principal components within the class of all random variables with at most $k (\leq p)$ dimensions. In this formulation it is desired to approximate the p component vector \underline{X} by a linear form $A \underline{Y}$ of a $k \times 1$ random vector \underline{Y} where A is a $p \times k$ matrix of constants. The error of approximation

F , may be measured by the trace of the

$$\text{residual covariance matrix} = E(\underline{X} - A\underline{Y})(\underline{X} - A\underline{Y})' . \quad (1.1)$$

Darroch showed that F is minimized with respect to A and \underline{Y} when and only when

$$A\underline{Y} = \underline{P}_1(\underline{P}_1'\underline{X}) + \cdots + \underline{P}_k(\underline{P}_k'\underline{X}) \quad (1.2)$$

which is a linear transformation of a set of k first principal components. Then, the minimum value is

$$F = \lambda_{k+1} + \cdots + \lambda_p$$

where, $\underline{P}_1, \dots, \underline{P}_k$ is any set of k orthonormal eigenvectors of Σ corresponding, respectively, to the eigenvalues $\lambda_1, \dots, \lambda_k$.

Okamoto and Kanazawa (1968) generalized Darroch's result, by allowing F to be any function of the eigenvalues of the residual covariance matrix, which is strictly increasing in each of its p arguments. Examples of such functions are the trace and the Euclidean norm. They showed that F is minimized when and only when $A\underline{Y}$ is as in (1.2), and then $F = F(\lambda_{k+1}, \dots, \lambda_p, 0, \dots, 0)$. A nice review of these results appears in Okamoto (1969).

Our extensions concern the complete residual covariance matrix rather than just increasing functions of the eigenvalues. These latter results require invariance under orthogonal transformations and do not allow an investigator to single out special components of the observation vector. After first establishing some preliminary results in Section 2, we show in Section 3 that if the criterion for the fit is the non-negative definite partial ordering on the residual covariance matrix, then in general no overall optimal $A\underline{Y}$ exists. However Theorem 3.2 establishes that an optimal class does exist. The implications of this extension are discussed in Section 3.3.

Principal components may also be introduced sequentially one at a time, by setting $k = 1$ above and demanding that at the j th stage, $A_j \underline{Y}_j$, be uncorrelated with those selected at a previous stage and that the residual covariance matrix $E(\underline{X} - A_j \underline{Y}_j)(\underline{X} - A_j \underline{Y}_j)'$ have eigenvalues that are as small as possible. This is the approach taken in Section 4 leading to Theorem 4.1.

Section 5 illustrates how the results of Darroch (1965) and Okamoto and Kanazawa (1968) can be viewed as the population analogues to a generalization of Pearson's (1901) approach. Both Darroch (1965) and Okamoto and Kanazawa seem to have been unaware of this fact. We also obtain a result, Theorem 5.4, on the approximation of cross products matrices.

Lastly in Section 5.2 a seemingly new interpretation is given of the sample principal components in R^n .

2. A Preliminary Formulation

We are interested in approximating the vector \underline{X} by a random vector \underline{AY} when goodness is measured by the covariance (1.1). Different choices of \underline{AY} produce different residual covariance matrices and in order to compare them, we give the following definition.

A partial ordering in the class of all non-negative definite matrices A is defined by the relation $\underline{\geq}$ where

$$A \underline{\geq} B \text{ iff } A - B \in A. \quad (2.1)$$

The problem then becomes that of finding the \underline{AY} which makes (1.1) as small as possible.

We begin by noticing that there is no loss in assuming $E\underline{X} = \underline{0}$. If this is not the case and $E\underline{X} = \underline{\mu}$ say, then similar to the regression situation where the best fitting polynomial is not forced to pass through the origin, we consider the residual to be defined by

$$E(\underline{X} - \underline{\eta} - \underline{AY})(\underline{X} - \underline{\eta} - \underline{AY})' \quad (2.2)$$

where $\underline{\eta}$ is a $p \times 1$ unknown constant vector that one is allowed to vary when searching for a minimum.

If $E\underline{Y} \neq \underline{0}$, the term \underline{AEY} may be absorbed into the $\underline{\eta}$ by replacing $\underline{X} - \underline{\eta} - \underline{AY}$ by $\underline{X} - (\underline{\eta} + \underline{AEY}) - [\underline{A}(\underline{Y} - E\underline{Y})]$. Since $\underline{\eta}$ is arbitrary, $\underline{\eta} + \underline{AEY}$ is also arbitrary and consequently we may restrict attention to \underline{Y} variables with zero expectation. Then

$$\begin{aligned} E(\underline{X} - \underline{\eta} - \underline{AY})(\underline{X} - \underline{\eta} - \underline{AY})' &= E(\underline{X} - \underline{\mu} - \underline{AY} - (\underline{\eta} - \underline{\mu}))(\underline{X} - \underline{\mu} - \underline{AY} - (\underline{\eta} - \underline{\mu}))' \\ &= E(\underline{X} - \underline{\mu} - \underline{AY})(\underline{X} - \underline{\mu} - \underline{AY})' + (\underline{\eta} - \underline{\mu})(\underline{\eta} - \underline{\mu})' \\ &\geq E(\underline{X} - \underline{\mu} - \underline{AY})(\underline{X} - \underline{\mu} - \underline{AY})' \end{aligned} \quad (2.3)$$

with strict inequality unless $\underline{\eta} = \underline{\mu}$.

Thus, in (2.3), we must take $\eta = \mu$. By assuming \underline{X} has already been corrected for its mean we get

$$E(\underline{X} - \underline{AY})(\underline{X} - \underline{AY})' \quad (2.4)$$

which is just the matrix in the original definition (1.1). Consequently, we may assume

$$E\underline{X} = E\underline{Y} = \underline{0}.$$

Also without loss of generality we may assume $E(\underline{Y}\underline{Y}') = I_k$, for even if rank of $E(\underline{Y}\underline{Y}') < k$ we may replace \underline{AY} by $A^*\underline{Y}^*$ where \underline{Y}^* is such that $E(\underline{Y}^*\underline{Y}^{*'}) = I_k$ and $\underline{AY} = A^*\underline{Y}^*$ (a.s.).

Finally notice that if $r = \text{rank of } \Sigma$ and $r \leq k$, the problem is trivial since we obtain a perfect fit with $A = [\underline{P}_1, \dots, \underline{P}_r, \underline{0}, \dots, \underline{0}]$ and $\underline{Y}' = (\underline{P}_1'\underline{X}, \dots, \underline{P}_r'\underline{X}, \dots, \underline{P}_k'\underline{X})$ so $\underline{AY} = \underline{X}$ and

$$0 = E(\underline{X} - \underline{AY})(\underline{X} - \underline{AY})' \leq E(\underline{X} - \bar{\underline{AY}})(\underline{X} - \bar{\underline{AY}})' \quad (2.5)$$

with strict inequality for any choice of $\bar{\underline{AY}}$ unless $\underline{AY} = \bar{\underline{AY}}$ (a.s.).

Thus in the rest of this paper, unless otherwise stated, we will assume that

$$E\underline{X} = E\underline{Y} = \underline{0}, \quad E\underline{Y}\underline{Y}' = I_k, \quad k < r. \quad (2.6)$$

Lemma 2.1. Under assumptions (2.6) let

$$B = \text{cov}(\underline{X}, \underline{Y}). \quad (2.7)$$

Then with strict inequality unless $B = A$

$$E(\underline{X} - \underline{AY})(\underline{X} - \underline{AY})' \geq \Sigma - BB' \geq 0. \quad (2.8)$$

Proof.

$$E(\underline{X} - \underline{AY})(\underline{X} - \underline{AY})' = E(\underline{X} - B\underline{Y} + (B - A)\underline{Y})(\underline{X} - B\underline{Y} + (B - A)\underline{Y})' \geq E(\underline{X} - B\underline{Y})(\underline{X} - B\underline{Y})' = \Sigma - BB'. \quad \square$$

3. A Class of Best Approximators

The best current result, on the approximation of \underline{X} by $B\underline{Y}$, is due to Okamoto and Kanazawa (1968). However, they restrict themselves to functions of eigenvalues of the residual matrix $E(\underline{X} - B\underline{Y})(\underline{X} - B\underline{Y})'$ which are increasing in each argument.

Theorem 3.1 [Okamoto and Kanazawa]. Let Σ have eigenvectors P_1, \dots, P_p corresponding to the eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Under assumptions (2.6), any strictly increasing function of the eigenvalues of

$$E(\underline{X} - A\underline{Y})(\underline{X} - A\underline{Y})'$$

is minimized with respect to A and \underline{Y} by the choice

$$A\underline{Y} = P_1 P_1' \underline{X} + \dots + P_k P_k' \underline{X}.$$

That is, the eigenvalues of the residual matrix are simultaneously minimized. \square

Remark. One can conclude from Theorem 3.1 that the sum of residual variances (trace) is minimized, the generalized variance (det.) is minimized, or the sum of squares of all entries (sum of squared eigenvalues) is minimized.

We show below that dropping the invariance condition, implied by the restriction to functions of eigenvalues, leads to a whole class of optimal solutions derived from the partial ordering of non-negative definiteness.

3.1. Best Approximators

Lemma 2.1 tells us that, in trying to minimize (1.1) under assumptions (2.6), we can only improve the approximation if the matrix of A coefficients of \underline{Y} is chosen as the covariance between \underline{X} and \underline{Y} . Now B must be found such that $\Sigma - BB'$ is as "small" as possible and then a \underline{Y} that produces such a B found. The next theorem will give us the structure of the non-negative definite covariance matrix $\Sigma - BB'$. It will be seen that there is no B which uniformly minimizes (2.8) but rather a collection of B 's which are optimal in a sense.

Theorem 3.2. Let $\Sigma \geq 0$ be of order p and rank r and have eigenvectors $P = [P_1, \dots, P_p]$ and eigenvalues $\lambda_1 \geq \dots \geq \lambda_p \geq 0$. Then under assumptions (2.6), set $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_r)$ and

$$\Gamma_k = \{B^* \underline{Y}^* | B^* \underline{Y}^* = P \begin{bmatrix} \Lambda^{1/2} R \begin{bmatrix} I_k & 0 \\ 0 & 0 \end{bmatrix} R' \Lambda^{-1/2} & 0 \\ 0 & 0 \end{bmatrix} P' \underline{X}; R \text{ } r \times r \text{ orthogonal}\}$$

so that by changing $R \begin{bmatrix} I_k & 0 \\ 0 & 0 \end{bmatrix} R'$ the different elements of Γ_k are obtained. Then

(1) [Completeness of Γ_k]. Given any random vector AY , there exists an element B^*Y^* in the class such that

$$E(\underline{X} - B^*Y^*)(\underline{X} - B^*Y^*)' \leq E(\underline{X} - AY)(\underline{X} - AY)'$$

with strict inequality unless (a.s.) $AY = B^*Y^*$.

(2) [Minimality of Γ_k]. Let $B\bar{Y}$, in Γ_k , be given. There is no $\bar{B}\bar{Y}$ in Γ_k (with $\bar{B}\bar{Y}$ not a.s. equal to B^*Y^*) such that

$$E(\underline{X} - \bar{B}\bar{Y})(\underline{X} - \bar{B}\bar{Y})' \leq E(\underline{X} - B^*Y^*)(\underline{X} - B^*Y^*)' . \quad \square$$

3.2. Proof of Theorem 3.2

We now proceed to develop a proof of Theorem 3.2 through a series of results.

Theorem 3.3. Let Σ of order p and rank r be n.n.d. and B be $p \times k$ such that $\Sigma - BB' \geq 0$. Then there exists an orthogonal matrix R of order r such that

$$BB' = P \begin{bmatrix} \Lambda^{1/2} R \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix} R' \Lambda^{1/2} & 0 \\ 0 & 0 \end{bmatrix} P'$$

where $D = \text{diag}(d_1, \dots, d_k)$; $1 \geq d_1 \geq \dots \geq d_k \geq 0$; $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_r)$ and P is orthogonal of order p and its i th column P_i is an eigenvector of Σ , corresponding to the i th largest eigenvalue of Σ , namely λ_i ($i = 1, \dots, p$).

Proof. By the spectral decomposition theorem, an orthogonal matrix P exists such that

$$\Sigma = P \begin{bmatrix} \Lambda & 0 \\ 0 & 0 \end{bmatrix} P'$$

where Λ is a $r \times r$ diagonal matrix with diagonal entries $\lambda_1 \geq \dots \geq \lambda_r \geq 0$. Then

$$\Sigma - BB' = P \begin{bmatrix} \Lambda & 0 \\ 0 & 0 \end{bmatrix} P' - BB' \geq 0 \iff \begin{bmatrix} \Lambda & 0 \\ 0 & 0 \end{bmatrix} - P'BB'P \geq 0 . \quad (3.1)$$

From (3.1) it is clear that the lower right hand corner of $P'BB'P$ must be zero, that is

$$P'B = \begin{bmatrix} C \\ 0 \end{bmatrix} \quad (3.2)$$

for some $r \times k$ matrix C . Equation (3.2) then takes the form

$$\begin{bmatrix} \Lambda & 0 \\ 0 & 0 \end{bmatrix} - \begin{bmatrix} CC' & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} \Lambda - CC' & 0 \\ 0 & 0 \end{bmatrix} \quad (3.3)$$

and we may thus restrict attention to the upper left corners.

$$\Lambda - CC' \geq 0 \iff I_r - \Lambda^{-1/2} CC' \Lambda^{-1/2} \geq 0. \quad (3.4)$$

Again by the spectral theorem an orthogonal matrix R exists such that

$$R' \Lambda^{-1/2} CC' \Lambda^{-1/2} R = \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix} \quad (3.5)$$

where $D = \text{diag}(d_1, \dots, d_k)$; $d_1 \geq \dots \geq d_k$. Substituting in (3.4)

$$I_r - \Lambda^{-1/2} CC' \Lambda^{-1/2} \geq 0 \iff I_r - \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix} \geq 0. \quad (3.6)$$

From (3.2)

$$B = P \begin{bmatrix} C \\ 0 \end{bmatrix}$$

so

$$BB' = P \begin{bmatrix} CC' & 0 \\ 0 & 0 \end{bmatrix} P'. \quad (3.7)$$

Equation (3.5) gives

$$CC' = \Lambda^{1/2} R \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix} R' \Lambda^{1/2} \quad (3.8)$$

which in (3.7) yields

$$BB' = P \begin{bmatrix} \Lambda^{1/2} R \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix} R' \Lambda^{1/2} & 0 \\ 0 & 0 \end{bmatrix} P'.$$

Finally, (3.6) shows $1 \geq d_1 \geq \dots \geq d_k \geq 0$. \square

Theorem 3.4 shows that the choice $D = I_k$ is best.

Theorem 3.4. Under the conditions of Theorem 3.3, let

$$B^*B^{*'} = P \begin{bmatrix} \Lambda^{1/2} R \begin{bmatrix} I_k & 0 \\ 0 & 0 \end{bmatrix} R' \Lambda^{1/2} & 0 \\ 0 & 0 \end{bmatrix} P' . \quad (3.9)$$

Then

$$(i) \quad 0 \leq \Sigma - B^*B^{*'} \leq \Sigma - BB' \quad (3.10)$$

with strict inequality unless $B^*B^{*'} = BB'$.

(ii) (Admissibility of B^*). If \bar{B} is an arbitrary $p \times k$ matrix such that

$$0 \leq \Sigma - \bar{B}\bar{B}' \leq \Sigma - B^*B^{*'} \quad (3.11)$$

then

$$\bar{B}\bar{B}' = B^*B^{*'} . \quad (3.12)$$

Proof. $(\Sigma - BB') - (\Sigma - B^*B^{*'}) = B^*B^{*'} - BB'$.

By Theorem 3.3 and (3.9), this can be written as

$$\begin{aligned} & P \begin{bmatrix} \Lambda^{1/2} R \begin{bmatrix} I_k & 0 \\ 0 & 0 \end{bmatrix} R' \Lambda^{1/2} & 0 \\ 0 & 0 \end{bmatrix} P' - P \begin{bmatrix} \Lambda^{1/2} R \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix} R' \Lambda^{1/2} & 0 \\ 0 & 0 \end{bmatrix} P' \\ & = P \begin{bmatrix} \Lambda^{1/2} R \begin{bmatrix} I_k - D & 0 \\ 0 & 0 \end{bmatrix} R' \Lambda^{1/2} & 0 \\ 0 & 0 \end{bmatrix} P' \geq 0 \end{aligned} \quad (3.13)$$

since, by (3.6), $1 \geq d_1 \geq \dots \geq d_k \geq 0$. Notice the inequality is strict unless $D = I_k$ in which case $BB' = B^*B^{*'}$.

Suppose \bar{B} exists such that (3.11) holds. Then, proceeding as in Theorem 3.3,

\bar{R} orthogonal and \bar{D} diagonal exist (the notation being clear) such that

$$\bar{B}\bar{B}' = P \begin{bmatrix} \Lambda^{1/2} \bar{R} \begin{bmatrix} \bar{D} & 0 \\ 0 & 0 \end{bmatrix} \bar{R}' \Lambda^{1/2} & 0 \\ 0 & 0 \end{bmatrix} P' . \quad (3.14)$$

From (i) of the present theorem

$$\bar{B}^* \bar{B}^{*'} = P \begin{bmatrix} \Lambda^{1/2} \bar{R} \begin{bmatrix} I_k & 0 \\ 0 & 0 \end{bmatrix} \bar{R}' \Lambda^{1/2} & 0 \\ 0 & 0 \end{bmatrix} P' \quad (3.15)$$

is such that

$$0 \leq \Sigma - \bar{B}^* \bar{B}^{*'} \leq \Sigma - \bar{B} \bar{B}' \leq \Sigma - B^* B^{*'} . \quad (3.16)$$

Equations (3.9) and (3.15) show that (3.16) is equivalent to

$$\bar{R} \begin{bmatrix} I_k & 0 \\ 0 & 0 \end{bmatrix} \bar{R}' \geq R \begin{bmatrix} I_k & 0 \\ 0 & 0 \end{bmatrix} R' . \quad (3.17)$$

It is easy to see that (3.17) can hold only if the equality holds, for it is clear that the first column of R , say r_1 , is an eigenvector of the right hand side of (3.17) corresponding to the eigenvalue 1. Since the eigenvalues are 1's and 0's, by the Courant-Fischer minimax theorem,

$$r_1' r_1 = 1 \leq r_1' \bar{R} \begin{bmatrix} I_k & 0 \\ 0 & 0 \end{bmatrix} \bar{R}' r_1 \leq 1 .$$

For equality, r_1 must be an eigenvector of the left hand side of (3.17) corresponding to the eigenvalue 1. Similarly it can be shown that the other columns of R are eigenvectors of both sides corresponding to the same eigenvalues and thus the two sides have the same eigenvalues and a common set of orthonormal eigenvectors which implies

$$\bar{R} \begin{bmatrix} I_k & 0 \\ 0 & 0 \end{bmatrix} \bar{R}' = R \begin{bmatrix} I_k & 0 \\ 0 & 0 \end{bmatrix} R' . \quad (3.18)$$

(Notice however this does not imply $R = \bar{R}$). \square

The next lemma gives us the general form of B^* .

Lemma 3.5. If $B^* B^{*'}$ is as in (3.9) and B^* is $p \times k$ then B^* is of the form

$$B^* = P \begin{bmatrix} \Lambda^{1/2} R \begin{bmatrix} Q \\ 0 \end{bmatrix} \\ 0 \end{bmatrix} \quad (3.19)$$

where Q is an arbitrary orthogonal matrix of order k .

Proof. We can rewrite (3.9) as

$$B^* B^{*'} = P \begin{bmatrix} \Lambda^{1/2} R \begin{bmatrix} I_k \\ 0 \end{bmatrix} \\ 0 \end{bmatrix} \left[\begin{bmatrix} I_k & 0 \end{bmatrix} R' \Lambda^{1/2}, 0 \right] P' . \quad (3.20)$$

Vinograd's theorem then implies that

$$B^* = P \begin{bmatrix} \Lambda^{1/2} R \begin{bmatrix} I_k \\ 0 \end{bmatrix} \\ 0 \end{bmatrix} Q = P \begin{bmatrix} \Lambda^{1/2} R \begin{bmatrix} Q \\ 0 \end{bmatrix} \\ 0 \end{bmatrix} \quad (3.21)$$

where Q is $k \times k$ orthogonal. \square

It remains to find \underline{Y}^* which satisfies the conditions in (2.6) and such that $E(\underline{X} \underline{Y}^{*'}) = B^*$. Although the argument that follows is very similar to that of Darroch (1965), it is given for completeness.

Lemma 3.6. Let

$$H = P \begin{bmatrix} \Lambda^{-1/2} R \begin{bmatrix} Q \\ 0 \end{bmatrix} \\ 0 \end{bmatrix} \quad (3.22)$$

and B^* be as in (3.21). Define $\underline{Y}^* = H' \underline{X}$. Then \underline{Y}^* is (a.s.) the only $k \times 1$ random vector satisfying

$$\begin{aligned} E(\underline{Y}^*) &= 0 \\ E(\underline{Y}^* \underline{Y}^{*'}) &= I_k \\ E(\underline{X} \underline{Y}^{*'}) &= B^* . \end{aligned} \quad (3.23)$$

Proof. Since

$$\Sigma H = P \begin{bmatrix} \Lambda & 0 \\ 0 & 0 \end{bmatrix} P' \cdot P \begin{bmatrix} \Lambda^{-1/2} R \begin{bmatrix} Q \\ 0 \end{bmatrix} \\ 0 \end{bmatrix} = P \begin{bmatrix} \Lambda^{1/2} R \begin{bmatrix} Q \\ 0 \end{bmatrix} \\ 0 \end{bmatrix} = B^*$$

and

$$H'B^* = [(Q' : 0)R'\Lambda^{-1/2}, 0]P'P \begin{bmatrix} \Lambda^{1/2}R \begin{bmatrix} Q \\ 0 \end{bmatrix} \\ 0 \end{bmatrix} = I_k \quad (3.24)$$

the conditions on the moments follow.

If \underline{Y} also satisfies (3.23), then

$$E(\underline{Y} - H'\underline{X})(\underline{Y} - H'\underline{X})' = I_k - E(\underline{Y}\underline{X}'H) - E(H'\underline{X}\underline{Y}') + I_k = 0 \quad (3.25)$$

so $\underline{Y} = H'\underline{X} = \underline{Y}^*$ (a.s.).

Remark. It must be noticed that even though (3.21) shows that B^* is arbitrary up to a choice of Q so \underline{Y}^* in Lemma 3.6 may change, the product $B^*\underline{Y}^*$ is invariant.

$$\begin{aligned} B^*\underline{Y}^* &= P \begin{bmatrix} \Lambda^{1/2}R \begin{bmatrix} Q \\ 0 \end{bmatrix} \\ 0 \end{bmatrix} [(Q' : 0)R'\Lambda^{-1/2} 0]P'\underline{X} \\ &= P \begin{bmatrix} \Lambda^{1/2}R \begin{bmatrix} I_k & 0 \\ 0 & 0 \end{bmatrix} R'\Lambda^{-1/2} & 0 \\ 0 & 0 \end{bmatrix} P'\underline{X} \end{aligned} \quad (3.26)$$

which does not depend on Q . We are now ready to prove our main result.

Proof (Theorem 3.2). (1) Let $A\underline{Y}$ be given. With $B = \text{cov}(\underline{X}, \underline{Y})$, Lemma 2.1 establishes that

$$E(\underline{X} - B\underline{Y})(\underline{X} - B\underline{Y})' = \Sigma - BB' \leq E(\underline{X} - A\underline{Y})(\underline{X} - A\underline{Y})' \quad (3.27)$$

with strict inequality unless $B = A$.

Since $\Sigma - BB' \geq 0$ from (2.8), Theorem 3.4 gives B^* such that

$$0 \leq \Sigma - B^*B^{*'} \leq \Sigma - BB' \quad (3.28)$$

with strict inequality unless $B^*B^{*'} = BB'$. Next, let \underline{Y}^* be as in Lemma 3.6 so

$$E(\underline{X} - B^*\underline{Y}^*)(\underline{X} - B^*\underline{Y}^*)' = \Sigma - B^*B^{*'} \leq E(\underline{X} - A\underline{Y})(\underline{X} - A\underline{Y})' \quad (3.29)$$

with strict inequality unless $B^*B^{*'} = BB'$ and $B = A$. Suppose that equality holds.

Then $B^*B^{*'} = BB'$ and by Lemma 3.5 applied to BB' , there exists an orthogonal Q , of order k , such that

$$B = P \begin{bmatrix} \Lambda^{1/2} R \begin{bmatrix} Q \\ 0 \end{bmatrix} \\ 0 \end{bmatrix} \quad (3.30)$$

Lemma 3.6 then implies that (a.s.)

$$\underline{Y} = ((Q' : 0) R' \Lambda^{-1/2}, 0) P' \underline{X} \quad (3.31)$$

and by the remark after Lemma 3.6, $B^* \underline{Y}^* = B \underline{Y}$ (a.s.).

For the proof of (2), we construct $\bar{B}^* \bar{Y}^*$, just as $B^* \underline{Y}^*$ was constructed in part (1), such that

$$E(\underline{X} - \bar{B}^* \bar{Y}^*)(\underline{X} - \bar{B}^* \bar{Y}^*)' = \Sigma - \bar{B}^* \bar{B}^{*'} \leq E(\underline{X} - B^* \underline{Y}^*)(\underline{X} - B^* \underline{Y}^*)' \quad (3.32)$$

That is, $\Sigma - \bar{B}^* \bar{B}^{*'} \leq \Sigma - B^* B^{*'}$ but then part (ii) of Theorem 3.4 implies that $\bar{B}^* \bar{B}^{*'} = B^* B^{*'}$. As in part (i), it then follows that $\bar{B}^* \bar{Y}^* = B^* \underline{Y}^*$ (a.s.).

3.3. Some Implications of Theorem 3.2

When predicting \underline{X} by $B \underline{Y}$ with \underline{Y} $k \times 1$, if loss is measured by a function of the eigenvalues of the residual covariance matrix, then Okamoto and Kanazawa (1968) give conditions under which principal components are optimal. Other objectives, which are not expressible as functions of the eigenvalues, lead to the selection of other members of the class Γ_k .

To illustrate this point, we consider the case $k = 1$. Suppose that, for a given vector \underline{a} , it is desired to find $B \underline{Y}$ such that

$$\underline{a}' E(\underline{X} - B \underline{Y})(\underline{X} - B \underline{Y})' \underline{a} \quad (3.33)$$

is a minimum. This monotone, non-decreasing loss function over the class of non-negative definite matrices is not a function of the eigenvalues. Here

$$\underline{a}' E(\underline{X} - B \underline{Y})(\underline{X} - B \underline{Y})' \underline{a} = \underline{a}' (\Sigma - B B') \underline{a} =$$

$$\text{Var}(\underline{a}' \underline{X}) - \text{Cov}^2(\underline{a}' \underline{X}, \underline{Y}) \geq \text{Var}(\underline{a}' \underline{X}) - \text{Var}(\underline{a}' \underline{X}) \text{Var}(\underline{Y}) = 0$$

with equality if and only if $\underline{Y} \propto \underline{a}' \underline{X}$. That is if $\underline{Y} = (\underline{a}' \Sigma \underline{a})^{-1/2} \underline{a}' \underline{X}$. Then

$$B = E(\underline{X} \underline{X}' \underline{a} / (\underline{a}' \Sigma \underline{a})^{1/2}) = (\underline{a}' \Sigma \underline{a})^{-1/2} \Sigma \underline{a}$$

$$B \underline{Y} = (\underline{a}' \Sigma \underline{a})^{-1} (\underline{a}' \underline{X}) \Sigma \underline{a} \quad (3.34)$$

with this choice of BY , the residual covariance matrix is

$$\Sigma - BB' = \Sigma - (a' \Sigma a)^{-1} \Sigma a a' \Sigma \quad (3.35)$$

and the loss in (3.33) is zero.

Now let PAP' be the spectral decomposition of Σ (we assume $|\Sigma| \neq 0$) with P orthogonal and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$; $\lambda_1 \geq \dots \geq \lambda_p$. We may write $a = cP\Lambda^{-1/2}r$ where c is a constant and $r \in R^p$, $r'r = 1$. Then, BY may be expressed as

$$BY = (a' \Sigma a)^{-1} (a' X) \Sigma a = P\Lambda^{-1/2} r r' \Lambda^{-1/2} P' X = P\Lambda^{1/2} R \begin{bmatrix} 1 & 0' \\ 0 & 0 \end{bmatrix} R' \Lambda^{-1/2} P' X \quad (3.36)$$

where $R = (r, r_2, \dots, r_p)$ is orthogonal. Consequently BY is in the class Γ_1 defined in Theorem 3.2. That is, over all possible predictors, the loss (3.33) is minimized when BY is as in (3.34) which is a member of Γ_1 .

As a further specialization, suppose our main concern is variable one. To reflect this, we may take $a' = (1, 0, \dots, 0)$ so that (3.33) equals entry (1,1) of the residual covariance matrix. In this situation, we may simply take $(X_1, 0, \dots, 0)$ as our predictor and have $X - (X_1, \dots, 0)' = (0, X_2, \dots, X_p)'$, which has residual covariance matrix with entry (1,1) equal to zero. However the predictor $(X_1, 0, \dots, 0)'$ does nothing to predict the remaining variables and conceivably, one can find a predictor that is as good at predicting X_1 and yet gives better prediction of the other variables. Any other predictor will give a residual covariance matrix with entry (1,1) greater than zero, unless the predictor's first entry is (a.s.) X_1 . We are then limited to predictors of the form AX_1 with A a $p \times 1$ vector having 1 as the first entry. From (3.34), we obtain the predictor

$$\frac{1}{\sigma_{11}} \Sigma_{(1)} X_1 \quad (3.37)$$

where $\Sigma_{(1)}$ is the first column of Σ , and σ_{11} its (1,1) entry. This predictor also gives a residual covariance matrix with entry (1,1) equal to zero.

We conclude our discussion by showing that, among all predictors AX_1 with first entry X_1 , (3.37) gives a residual covariance matrix whose entries are smaller than

or equal to those corresponding to the residual covariance matrix given by any other predictor of the prescribed form.

To this end let a predictor AX_1 be written as

$$\underline{X} = \begin{bmatrix} X_1 \\ Z \end{bmatrix}, \quad AX_1 = \begin{bmatrix} X_1 \\ CX_1 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \sigma_{11} & \sigma_1' \\ \sigma_1 & \Sigma_{11} \end{bmatrix}$$

then $\underline{X} - AX_1$ has residual covariance matrix

$$\begin{bmatrix} 0 & 0' \\ 0 & \Sigma_{11} - \sigma_1 C' - C \sigma_1' + C \sigma_{11} C' \end{bmatrix}$$

but

$$\Sigma_{11} - \sigma_1 C' - C \sigma_1' + C \sigma_{11} C' = \left[\frac{\sigma_1}{\sqrt{\sigma_{11}}} - \sqrt{\sigma_{11}} C \right] \left[\frac{\sigma_1}{\sqrt{\sigma_{11}}} - \sqrt{\sigma_{11}} C \right]' + \Sigma_{11} - \frac{1}{\sigma_{11}} \sigma_1 \sigma_1'$$

$$\geq \Sigma_{11} - \frac{1}{\sigma_{11}} \sigma_1 \sigma_1'.$$

4. Introducing the Principal Components Sequentially

Here we extend Theorem 3.1 to the sequential selection of approximators. Suppose $k = 1$, so that the covariance matrices are of the form $E(\underline{X} - AY)(\underline{X} - AY)'$ where A is a $p \times 1$ vector of constants and Y a univariate random variable. In this section, we sometimes write $\lambda[M]$ to denote that λ is an eigenvalue of M .

Let P_1 be an eigenvector of Σ , corresponding to its largest eigenvalue, and $A_1 Y_1 = P_1 P_1' \underline{X}$. Theorem 3.1 establishes that this AY has the property that the i th largest eigenvalue of the corresponding residual covariance matrix $E(\underline{X} - P_1 P_1' \underline{X})(\underline{X} - P_1 P_1' \underline{X})'$ is less than or equal to the corresponding eigenvalue of any possible residual covariance matrix ($i = 1, \dots, p$).

Proceeding in a sequential manner, we seek $A_2 Y_2$ uncorrelated with $A_1 Y_1$ with the same minimizing property among the class of AY 's uncorrelated with $A_1 Y_1$. Then $A_3 Y_3$ uncorrelated with $A_1 Y_1$ and $A_2 Y_2$, and with the minimizing property is desired, and so on. The next theorem tells us that principal components are the solution.

Theorem 4.1. Let \underline{X} be a $p \times 1$ random vector with zero expectation and covariance matrix Σ of rank r . Also let $P = \{AY | A \text{ is a } p \times 1 \text{ vector and } Y \text{ a random variable with } EY = 0; EY^2 = 1\}$. For $j = 1, \dots, r$, let $A_j Y_j \in P$. Then,

$$\lambda_i [E(\underline{X} - A_j Y_j)(\underline{X} - A_j Y_j)'] \leq \lambda_i [E(\underline{X} - AY)(\underline{X} - AY)'] \quad (i = 1, \dots, p),$$

for every AY in P uncorrelated with $A_1 Y_1, \dots, A_{j-1} Y_{j-1}$, if and only if $A_j Y_j = P_j P_j' \underline{X}$ where P_j is an eigenvector of Σ corresponding to the eigenvalue $\lambda_j[\Sigma]$.

Proof. For $j = 1$, application of Theorem 3.1 provides $A_1 Y_1 = P_1 P_1' \underline{X}$ where P_1 is as stated. For $j = 2$. Let $P = (P_1 : C) = (P_1, P_2, \dots, P_p)$ be orthogonal where P_2, \dots, P_p are eigenvectors of Σ corresponding respectively to $\lambda_2[\Sigma], \dots, \lambda_p[\Sigma]$, so that $\underline{X} = PP' \underline{X} = P_1 P_1' \underline{X} + CC' \underline{X}$ and

$$\begin{aligned} E(\underline{X} - A_2 Y_2)(\underline{X} - A_2 Y_2)' &= E(P_1 P_1' \underline{X} + CC' \underline{X} - A_2 Y_2)(P_1 P_1' \underline{X} + CC' \underline{X} - A_2 Y_2)' \\ &= E(P_1 P_1' \underline{X})(P_1 P_1' \underline{X})' + E(CC' \underline{X} - A_2 Y_2)(CC' \underline{X} - A_2 Y_2)' \\ &= \lambda_1 P_1 P_1' + E(CC' \underline{X} - A_2 Y_2)(CC' \underline{X} - A_2 Y_2)' \end{aligned}$$

where the cross terms are zero because $P_1 P_1' \underline{X}$ is uncorrelated with $CC' \underline{X}$ and $A_2 Y_2$.

The result now follows from Theorem 3.1 by noting that $CC' \underline{X}$ has covariance matrix

$$CC' \Sigma CC' = CC' (P_1 : C) \Lambda \begin{bmatrix} P_1' \\ C' \end{bmatrix} CC' = \sum_{i=2}^p \lambda_i P_i P_i' = C \begin{pmatrix} \lambda_2 & & \\ & \ddots & \\ & & \lambda_p \end{pmatrix} C'.$$

Alternatively, if $B_2 = E(CC' \underline{X} Y_2) = CC' E(\underline{X} Y_2) = CC' u$, say. From Lemma 2.1, we have

$$\begin{aligned} \lambda_1 P_1 P_1' + E(CC' \underline{X} - A_2 Y_2)(CC' \underline{X} - A_2 Y_2)' &\geq \lambda_1 P_1 P_1' + CC' \Sigma CC' - B_2 B_2' \\ &= \lambda_1 P_1 P_1' + CC' (\Sigma - uu') CC'. \end{aligned} \quad (4.1)$$

Since $C' P_1 = 0$, by simply multiplying (4.1) on the right by P_1 we see that P_1 is an eigenvector corresponding to the eigenvalue λ_1 . Similarly any eigenvector of the second term, perpendicular to P_1 , is also an eigenvector of the whole right hand side of (4.1) for the same eigenvalue. Thus to minimize the eigenvalues of (4.1) we need only minimize the eigenvalues of the second term.

But

$$CC'(\Sigma - uu')CC' = E(CC'X - A_2Y_2)(CC'X - A_2Y_2)'$$

and $CC'X$ is a vector with zero expectation and covariance matrix $CC'\Sigma CC'$. From the result for $j = 1$, we must take $A_2Y_2 = u_2u_2'X$ where u_2 is an eigenvector of $CC'\Sigma CC'$. Thus $A_2Y_2 = P_2P_2'X$ gives the result for $j = 2$. Similarly for $j = 2, \dots, r$. \square

Remark. The procedure was only carried through r stages where $r = \text{rank of } \Sigma$, instead of p stages. The reason for this is that $X = P_1P_1'X + \dots + P_rP_r'X$ (a.s.) so that any variable uncorrelated with $P_1P_1'X, \dots, P_rP_r'X$ is uncorrelated with X and then

$$E(X - AY)(X - AY)' = \Sigma + AE(Y^2)A' \geq \Sigma.$$

If one insists on introducing p components, then since (a.s.)

$P_{r+1}'X = \dots = P_p'X = 0$, $P_{r+1}P_{r+1}'X, \dots, P_pP_p'X$ may be taken as the extra components. However no matter how the extra components are introduced, stages $r + 1, \dots, p$ are irrelevant.

Remark. An alternative way of sequentially introducing the first r ($= \text{rank of } \Sigma$) principal components is to begin as above with $A_1Y_1 = P_1P_1'X$ having the property that its eigenvalues are less than or equal to the corresponding eigenvalues of any possible residual covariance matrix. However, as a second step, we consider the residual

$X - P_1P_1'X = Z$ which itself has zero expectation and covariance matrix

$\Sigma - \lambda_1 P_1P_1' = \sum_{i=2}^p \lambda_i P_iP_i'$. Applying Theorem 3.1 to the residual Z , one finds that $AY = P_2P_2'X$ gives a corresponding residual covariance matrix

$$E(Z - P_2P_2'X)(Z - P_2P_2'X)' = \Sigma - \lambda_1 P_1P_1' - \lambda_2 P_2P_2'$$

such that for $i = 1, \dots, p$, its i th largest eigenvalue is less or equal than the corresponding eigenvalue of any other possible residual covariance matrix. In this way, the first r principal components may be introduced sequentially in terms of approximating successive residuals. We then have

Corollary 4.2. Suppose we select the approximators one at a time from the residual covariance of the previous stage. Let $Z_q = X - \sum_{i=1}^q P_iP_i'X$ be the residual matrix after q steps. Then the choice $AY = P_{q+1}P_{q+1}'X$, for stage $q + 1$, minimizes the eigenvalues of $E(Z_q - AY)(Z_q - AY)'$.

5. Some Sample Interpretations

5.1. p-Space Interpretations

Let $(x_1, \dots, x_n) = x$ denote n observations on the $p \times 1$ random vector X . We can think of the n observations as n points in R^p . Moreover, one can assign probability $1/n$ to each observation vector and apply the results from the population model to obtain the sample results. We intend, however, to go in the reverse direction and show how Pearson's (1901) result may be extended to hold for all increasing functions of eigenvalues. This frames the result by Okamoto and Kanazawa so that it may be viewed as a natural population analog. In our development, although we speak of projections on planes not passing through the origin we really make suitable translations and take proper projections on subspaces.

Suppose that we are interested in finding that line which best fits the n given points, in the sense that the sum of the squares of the distances from these n points to the line is a minimum. More generally, we ask for the hyperplane of dimension $k (\leq p)$ which best fits the data in the above sum of squares sense.

Let $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$, be the centroid of the n points and $S = \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})'$ the cross products matrix. We have the following theorem due to Pearson (1901).

Theorem 5.1 [Pearson]. The hyperplane of dimension $k (\leq p)$ such that the sum of squares of the distances from the points to the plane is a minimum is of the form

$$\{x | x = \bar{x} + P\mu, \mu \in R^k\}$$

where $P = (P_1, \dots, P_k)$ and its k columns constitute a set of k orthonormal eigenvectors of the cross products matrix S corresponding respectively to the k largest eigenvalues. \square

If $\bar{x} + P\mu_i$ is the point in the plane closest to x_i then

$$P\mu_i = PP'(x_i - \bar{x}) = P_1 P_1'(x_i - \bar{x}) + \dots + P_k P_k'(x_i - \bar{x}) \quad (5.1)$$

and if $\lambda_1 \geq \dots \geq \lambda_p$ are the eigenvalues of S then the sum of squared distances from the points x_i to the hyperplane equals

$$SS(\text{distances}) = \lambda_{k+1} + \dots + \lambda_p. \quad (5.2)$$

Pearson's result then tells us that the "best" hyperplane is determined by \bar{x}, P_1, \dots, P_k and the point nearest to x_i is $\bar{x} + P_1 P_1' (x_i - \bar{x}) + \dots + P_k P_k' (x_i - \bar{x})$.

We may also think of P_i as determining a privileged direction in R^p . The line $\{P_i \alpha \mid \alpha \in R\}$ is the line perpendicular to P_1, \dots, P_{i-1} which best fits the points $(x_1 - \bar{x}, \dots, x_n - \bar{x}) = \chi^*$ or equivalently the line that best fits (no restrictions now) the residuals $\chi^* - P_1 P_1' \chi^* \dots - P_{i-1} P_{i-1}' \chi^*$. We now show how Pearson's result can be extended to a statement about eigenvalues.

Let P be a $p \times k$ matrix, $\mu = (\mu_1, \dots, \mu_n)$ be a $k \times n$ matrix and \bar{x} a $p \times 1$ vector, then the points $\bar{x} + P \mu_i$ ($i = 1, \dots, n$) all lie in the hyperplane parallel to that spanned by the columns of P . Consider the matrix

$$\sum_{i=1}^n (x_i - \bar{x} - P \mu_i) (x_i - \bar{x} - P \mu_i)' \quad (5.3)$$

If $\bar{x} + P \mu_i$ is the projection of x_i over the hyperplane then the trace of (5.3) gives us the sum of squared distances from the x_i 's to the hyperplane and Pearson's result tells us how to select \bar{x} and P . Here we show how to make the eigenvalues of (5.3) as small as possible by a correct choice of $\bar{x}, P, \mu_1, \dots, \mu_n$. It is clear that the columns of P can be taken to be orthonormal in the following optimization.

Theorem 5.2. To minimize the eigenvalues of the matrix

$$\sum_{i=1}^n (x_i - \bar{x} - P \mu_i) (x_i - \bar{x} - P \mu_i)' \quad (5.4)$$

\bar{x} may be selected as \bar{x} , and

$$P(\mu_1, \dots, \mu_n) = P P' \chi^*$$

where χ^* stands for the matrix χ with each of its rows corrected for the mean and P_1, \dots, P_k constitute a set of k orthonormal eigenvectors of the cross products matrix S corresponding respectively to the k largest eigenvalues $\lambda_1 \geq \dots \geq \lambda_k$. Thus, any increasing function of the eigenvalues is minimized by this choice.

Proof. Without loss of generality, we may assume that $P \bar{\mu} = \sum_{i=1}^n P \mu_i = 0$ for, if this is not the case, \bar{x} may be replaced by $\bar{x} + P \bar{\mu}$ and $P \mu_i$ by $P \mu_i - P \bar{\mu}$. Suppose then, $P \bar{\mu} = 0$. We have the non-negative definite ordering

$$\begin{aligned}
\sum_i (x_i - \bar{x} - p\mu_i)(x_i - \bar{x} - p\mu_i)' &= \sum_i (x_i - \bar{x} - p\mu_i + \bar{x} - \bar{x})(x_i - \bar{x} - p\mu_i + \bar{x} - \bar{x})' \\
&= \sum_i (x_i - \bar{x} - p\mu_i)(x_i - \bar{x} - p\mu_i)' + \sum_i (\bar{x} - \bar{x})(\bar{x} - \bar{x})' \\
&\geq \sum_i (x_i - \bar{x} - p\mu_i)(x_i - \bar{x} - p\mu_i)' \quad (5.5)
\end{aligned}$$

with strict inequality unless $\bar{x} = \bar{x}$ so that \bar{x} must be taken as the centroid of x_1, \dots, x_n .

Next we write $Y = (y_1, \dots, y_n) = (x_1 - \bar{x}, \dots, x_n - \bar{x})$ and obtain

$$\sum_{i=1}^n (y_i - p\mu_i)(y_i - p\mu_i)' = (Y - p\mu)(Y - p\mu)', \quad \mu = (\mu_1, \dots, \mu_r).$$

The above matrix has the same nonzero eigenvalues as the matrix

$$(Y - p\mu)'(Y - p\mu)$$

and we have

$$\begin{aligned}
(Y - p\mu)'(Y - p\mu) &= (Y - PP'Y + PP'Y - p\mu)'(Y - PP'Y + PP'Y - p\mu) \\
&= (Y - PP'Y)'(Y - PP'Y) + (PP'Y - p\mu)'(PP'Y - p\mu) \\
&\geq (Y - PP'Y)'(Y - PP'Y)
\end{aligned}$$

with strict inequality unless $PP'Y = p\mu$ or $P'Y = \mu$. With this choice we have

$$(Y - PP'Y)'(Y - PP'Y) = Y'(I - PP')Y.$$

We want to minimize the eigenvalues of this matrix or equivalently of the matrix

$$(I - PP')YY'(I - PP') = (I - PP')S(I - PP') = QQ'SQQ'$$

where Q is such that $(P; Q)$ is orthogonal. Next, consider the product

$$\begin{bmatrix} P' \\ Q' \end{bmatrix} QQ'SQQ' \begin{bmatrix} P & Q \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & Q'SQ \end{bmatrix} \quad (5.7)$$

which has the same eigenvalues as $QQ'SQQ'$. Its eigenvalues are k zeroes and those of the lower right corner of

$$\begin{bmatrix} P' \\ Q' \end{bmatrix} S \begin{bmatrix} P & Q \end{bmatrix} = \begin{bmatrix} P'SP & P'SQ \\ Q'SP & Q'SQ \end{bmatrix}.$$

The Poincare separation theorem (c.f. Bellman (1970), p. 117) shows that

$$\lambda_1 [Q'SQ] \geq \lambda_{k+1} [S], \dots, \lambda_{p-k} [Q'SQ] \geq \lambda_p [S]$$

with at least one strict inequality unless an orthogonal matrix of the form

$$\begin{bmatrix} R_1 & 0 \\ 0 & R_2 \end{bmatrix}$$

with R_1 of order k exists, such that

$$\begin{bmatrix} R_1' & 0 \\ 0 & R_2' \end{bmatrix} \begin{bmatrix} P' \\ Q' \end{bmatrix} S \begin{bmatrix} P & Q \end{bmatrix} \begin{bmatrix} R_1 & 0 \\ 0 & R_2 \end{bmatrix} = \text{diag}(\lambda_1, \dots, \lambda_p)$$

where $\lambda_1 \geq \dots \geq \lambda_p$ are the eigenvalues of S .

Then the columns of PR_1 must constitute a set of k orthonormal eigenvectors of S corresponding respectively to the k largest eigenvalues $\lambda_1 \geq \dots \geq \lambda_k$. Finally $PP' = PR_1R_1'P'$ gives the desired result. \square

We now state the sample version of the result by Okamoto and Kanazawa (1968) or Okamoto (1969). Our intention is to illustrate how this is connected to Pearson (1901).

Theorem 5.3. Let \tilde{X} denote the random variable which assumes each of the observed sample values x_i , $i = 1, 2, \dots, n$ with probability $\frac{1}{n}$ and \tilde{U} denote any corresponding $k \times 1$ vector. Then, the eigenvalues of

$$E[\tilde{X} - \bar{x} - A\tilde{U}][\tilde{X} - \bar{x} - A\tilde{U}]' = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x} - Au_i)(x_i - \bar{x} - Au_i)'$$

are simultaneously minimized over all A of dimension $p \times k$ and all sets of values $\{u_i\}$ for \tilde{U} by

$$Au_i = (P_1P_1' + \dots + P_kP_k')(x_i - \bar{x})$$

where P_1, \dots, P_k is a set of orthonormal eigenvectors of the matrix $\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})'$ corresponding to the k largest eigenvalues.

Proof. The result follows directly from Theorem 5.2 or as a special case of Theorem 3.1. \square

We also have another geometric interpretation in terms of approximating cross product matrices.

Let a hyperplane be $\{v | v = \bar{x} + P\mu; \mu \in R^k\}$ where $P = (P_1, \dots, P_k)$ is $p \times k$ with orthonormal columns. The projection of a given point x_i into this hyperplane is the point $v_i = \bar{x} + PP'(x_i - \bar{x})$. Let $S = \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})'$ denote the cross products matrix and let

$$S_v = \sum_{i=1}^n (v_i - \bar{v})(v_i - \bar{v})' = \sum_{i=1}^n PP'(x_i - \bar{x})(x_i - \bar{x})'PP' = PP'SPP' \quad (5.8)$$

denote the cross products matrix of the projected points. We want to find a hyperplane in R^p , passing through the centroid \bar{x} , and such that S_v is close to S .

Theorem 5.4. In order to minimize the eigenvalues of the discrepancy in sums of cross product matrices

$$S - S_v = S - PP'SPP',$$

over all k dimensional planes passing through \bar{x} , we must take the plane where the columns of P are a set of eigenvectors of S corresponding to the k largest eigenvalues.

Proof. Let Q be such that $(P; Q)$ is orthogonal. The eigenvalues of $S - PP'SPP'$ are exactly the same as those of

$$\begin{bmatrix} P' \\ Q' \end{bmatrix} S (P \ Q) - \begin{bmatrix} P' \\ Q' \end{bmatrix} PP'SPP' (P \ Q) = \begin{bmatrix} 0 & P'SQ \\ Q'SP & Q'SQ \end{bmatrix}$$

whose eigenvalues are k zeroes together with those of $Q'SQ$.

From the argument given after equation (5.7), it follows that the hyperplane has to be of the form

$$\{v | v = \bar{x} + Pu; u \in R^k\} \quad (5.9)$$

where $P = (P_1, \dots, P_k)$ and P_1, \dots, P_k constitute a set of orthonormal eigenvectors of the matrix S corresponding respectively to the k largest eigenvalues $\lambda_1 \geq \dots \geq \lambda_k$. \square

Remark. If, alternatively, we consider the residuals $x_i - v_i$ with v_i as in (5.8), then the cross products matrix for the residuals is

$$S_r = \sum (x_i - v_i)(x_i - v_i)' = \sum (x_i - \bar{x} - PP'(x_i - \bar{x}))(x_i - \bar{x} - PP'(x_i - \bar{x}))'$$

and Theorem 5.2 tells us that, to minimize the eigenvalues of S_r , the hyperplane must also be taken as in (5.9).

5.5. n-Space Interpretation

A geometrical interpretation of principal components in n space does not seem to have been given by previous workers. Let $x = (x_1, \dots, x_n)$ denote n observations on the p -vector \bar{x} . If the i -th row of $x = (x_1, \dots, x_n)$ is denoted by y_i' and the i -th row of $x^* = (x_1 - \bar{x}, \dots, x_n - \bar{x})$ by $y_i^{*'}$, then these rows are points in R^n .

With the mean corrected rows of the data matrix, it is natural to apply Euclidean distances and inner products. Suppose then, we ask for the plane of dimension $k (\leq p)$ which best fits the p points determined by the rows of χ^* and which passes through the origin. If

$$C = (c_1, \dots, c_n) = \begin{bmatrix} r'_1 \\ \vdots \\ r'_p \end{bmatrix}$$

denotes the points in the hyperplane with r'_i being the point closest to y_i^* , we want to minimize the sum of the squared norms of the rows of the matrix

$$\begin{bmatrix} y_1^* - r'_1 \\ \vdots \\ y_p^* - r'_p \end{bmatrix} = (\chi^* - C) = (x_1^* - c_1, \dots, x_n^* - c_n) \quad (5.10)$$

which is clearly equivalent to minimizing the sum of squares of all the entries of the above matrix or to minimizing the sum of the squared norms of the columns. From Pearson's result, the best choice of C is then

$$C = p_1 p_1' \chi^* + \dots + p_k p_k' \chi^* \quad (5.11)$$

where p_1, \dots, p_k are as in (5.1).

For the case $k = 1$, we can consider the problem as one of first projecting the rows of χ^* on any $n \times 1$ vector a' . These projections are given by the rows of

$$\chi^* a a' / (a' a) = C_a \quad (5.12)$$

where C_a is of rank 1 when $\text{rank}(\chi^*) \geq 1$. Thus, the sum of squares of the i -th row of $\chi^* - C_a$ is the squared distance from y_i^* to the line determined to a . Moreover, we know that the choice $a' = p_1' \chi^* = (p_1' (x_1 - \bar{x}), \dots, p_1' (x_n - \bar{x}))$, the sample values of the first principal component, produces C_a of the form (5.12) or

$$\chi^* \chi^{*'} p_1 p_1' \chi^* / p_1' \chi^* \chi^* p_1 = S p_1 p_1' \chi^* / p_1' S p_1 = \lambda_1 p_1 p_1' \chi^* / \lambda_1 = p_1 p_1' \chi^*$$

conforming to the optimal choice of C given by (5.11) with $k = 1$.

Remark. The first principal component minimizes the sum of squared distances to the rows of the mean corrected data matrix χ^* . This is illustrated in Figure 5.1. Our geometrical interpretation helps us understand the relative importance of a row with large length (sample variance) in determining the principal component.

Remark. If the rows of χ^* are first scaled to become unit vectors, then extracting principal components from the new matrix is equivalent to using the sample correlation matrix R . As the vectors in the geometric interpretation are now all of unit length, minimizing the sum of squared errors is equivalent to maximizing (minimizing) the sum of $\cos^2 \theta_i$ ($\sin^2 \theta_i$) where θ_i is the angle between y_i^* and \underline{a} .

It is easy to see that $P_1' \chi^*$ determines a privileged line in n space. It is the line perpendicular to $P_1' \chi^*, \dots, P_{i-1}' \chi^*$ that passes through the origin and which best fits the rows of χ^* or equivalently the line through the origin which best fits the rows (no restrictions now) of the residual $\chi^* - P_1 P_1' \chi^* \dots - P_{i-1} P_{i-1}' \chi^*$.

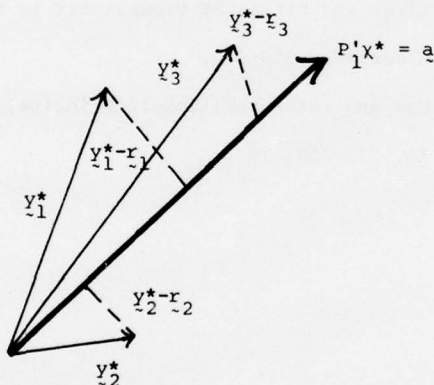


Figure 5.1. Case $p = 3$, showing sample first principal component minimizing the sum of squared distances

$$\|y_1^* - r_1\|^2 + \|y_2^* - r_2\|^2 + \|y_3^* - r_3\|^2$$

from the y_i^* 's to line.

Acknowledgment

The authors wish to thank Dr. M. El-Bassiouni for suggesting improvements in the presentation of this material.

REFERENCES

- Anderson, T. W. (1958). An Introduction to Multivariate Statistical Analysis. New York, Wiley.
- Bellman, R. (1970). An Introduction to Matrix Analysis. McGraw-Hill, New York.
- Darroch, J. N. (1965). An Optimal Property of Principal Components. Ann. Math. Statist., 36, 1579-1582.
- Girshick, M. A. (1936). Principal Components. J. Amer. Statist. Assoc., 31, 519-528.
- Hotelling, H. (1933). Analysis of Complex Statistical Variables into Principal Components. J. Educ. Psych., 24, 417-444, 498-520.
- Okamoto, M. (1969). Optimality of Principal Components, Multivariate Analysis II, 673-685, Ed. P. R. Krishnaiah, Academic Press, New York.
- Okamoto, M. and Kanazawa M. (1968). Minimization of Eigenvalues of a Matrix and Optimality of Principal Components. Ann. Math. Statist., 39, 859-863.
- Pearson, K. (1901). On Lines and Planes of Closest Fit to Systems of Points in Space. Phil. Mag., 2, (sixth series), 559-572.
- Rao, C. R. (1964). The Use and Interpretation of Principal Component Analysis in Applied Research. Sankhya, 26, 329-358.

14 MAR-TSR-1836

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

| REPORT DOCUMENTATION PAGE | | READ INSTRUCTIONS BEFORE COMPLETING FORM |
|---|--|---|
| 1. REPORT NUMBER 1836 | 2. GOVT ACCESSION NO. | 3. RECIPIENT'S CATALOG NUMBER |
| 4. TITLE (and Subtitle) SOME OPTIMAL PROPERTIES AND INTERPRETATIONS OF PRINCIPAL COMPONENTS. | 5. TYPE OF REPORT & PERIOD COVERED Summary Report - no specific reporting period | |
| 6. AUTHOR(s) Raul/Hudlet Richard A./Johnson | 7. CONTRACT OR GRANT NUMBER(s) DAAG29-75-C-0024, JNSF-MCS77-09574 | |
| 8. PERFORMING ORGANIZATION NAME AND ADDRESS Mathematics Research Center, University of 610 Walnut Street Wisconsin Madison, Wisconsin 53706 | 9. PROJECT ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS Work Unit # 4 - Probability, Statistics, and Combinatorics | |
| 10. CONTROLLING OFFICE NAME AND ADDRESS See Item 18 below. | 11. REPORT DATE March 1978 | |
| 12. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office) | 13. NUMBER OF PAGES 24 | |
| 14. DISTRIBUTION STATEMENT (of this Report) Approved for public release; distribution unlimited. | 15. SECURITY CLASS. (of this report) UNCLASSIFIED | |
| 15. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report) | 16a. DECLASSIFICATION/DOWNGRADING SCHEDULE | |
| 16. SUPPLEMENTARY NOTES U. S. Army Research Office National Science Foundation P. O. Box 12211 Washington, D. C. 20550 Research Triangle Park North Carolina 27709 | | |
| 17. KEY WORDS (Continue on reverse side if necessary and identify by block number) Principal components Statistical approximations | | |
| 18. ABSTRACT (Continue on reverse side if necessary and identify by block number) Optimal properties are derived and some new geometrical interpretations given for principal components. Typically, our main results concern the simultaneous minimization of eigenvalues of certain covariance matrices which measure the goodness of an approximation. Many popular criteria like total variance and generalized variance, which are increasing functions of the eigen- values, are then minimized by the best approximator. | | |

20. ABSTRACT - Cont'd.

→ In other situations, the criterion may not be a monotone function of the eigenvalues. ^{is derived} In ~~theorem 3.2~~, we ~~derive~~ a general optimal class based on the non-negative definite ordering of covariance matrices. ~~Theorem 4.1 gives~~ ^{is given} a result for the sequential selection of principal components. In the final section, ~~we give~~ a new geometrical interpretation of the sample principal components.

Also